

1 We are very grateful to the constructive comments. Now, we give responses for all questions/comments.

2 **To reviewer #1:** Thanks for your constructive comments to enhance the completeness of our work.

3 **Q1:** The comparison of validation error between search network and child network... The search progress...

4 **R1:** For different models, Table 1 in this response reports the comparison of the validation errors at the end of search
5 and after architecture derivation without fine-tuning. The results show that DATA (gap=1.87, err=9.21) is superior to
6 maintain better performance compared with SNAS (gap=2.13, err=9.33) and especially DARTS (gap=33.03, err=45.34),
7 while requiring less training epochs than SNAS (100 vs 150 epochs), as the search progresses illustrated in Figure 2.

8 **Q2:** It is complained that NAS approach has high variance over different initializations... A variance plot...

9 **R2:** In Table 2 and Figure 3, the stds of DATA and the variances of DATA with different sampling time (M) are reported,
10 which indicate that DATA (M=1, std=0.08) achieves lower stds than DARTS (k=1, std=0.14). Furthermore, Figure 4
11 shows the variances of architectures in search process for a comprehensive analysis. Please kindly check them.

12 **Q3:** Search result of M=4 and M=7.

13 **R3:** The search result (M=4) is shown in Figure 1. More results and analyses will be included in an updated version.

14 **Q4:** The proposed method is sampling based... And the gradient estimator should be written...

15 **R4:** Thanks for the suggestion, we will write the gradient estimator in an updated version. In brief, we have:

$$\begin{aligned}
 \frac{\partial f_{P_k^{(i,j)}(G^{(i,j)})}}{\partial P_k^{(i,j)}} &= \frac{\frac{\partial}{\partial P_k^{(i,j)}} \exp\left(\frac{\log P_k^{(i,j)} + G_k^{(i,j)}}{\tau}\right)}{\sum_{l=1}^K \exp\left(\frac{\log P_l^{(i,j)} + G_l^{(i,j)}}{\tau}\right)} \left(\delta(k' - k) - \frac{\exp\left(\frac{\log P_k^{(i,j)} + G_k^{(i,j)}}{\tau}\right)}{\sum_{l=1}^K \exp\left(\frac{\log P_l^{(i,j)} + G_l^{(i,j)}}{\tau}\right)} \right) \\
 &= \frac{\delta(k' - k) - f_{P_k^{(i,j)}(G^{(i,j)})}}{\tau P_k^{(i,j)}} f(G_k^{(i,j)}),
 \end{aligned}$$

18 where $G_k^{(i,j)}$ is the k -th Gumbel random variable, $P_k^{(i,j)}$ is the k -th element in $\mathbf{P}^{(i,j)} \in \mathbb{R}^K$ and it denotes the probability
19 of choosing the k -th operation on the edge $e^{(i,j)}$. Then, the gradient of the loss with respect to other parameters can be
20 easily calculated using the chain rule, since the other components in the network are differentiable as well.

21 **Q5:** Since the objective of architecture weights and probability vector are both training loss... is overfitting an issue...

22 **R5:** Similar to the dropout to combat overfitting, our EGS provides a network sampling technique that can be treated
23 as a specific regularization, which is in agreement with the motivation in [C1]. That is, EGS endows DATA with the
24 capability of searching architectures by learning on the sampled networks. Empirically, our experiments have strongly
25 evidenced its validity. For a clearer verification, the example you mentioned will be included in an updated version.

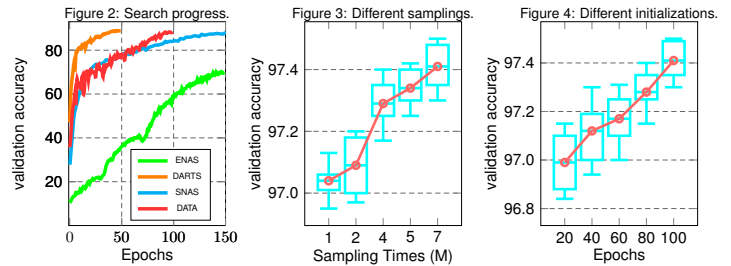
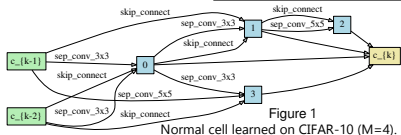
26 **Q6:** DATA should be slower in terms per search iteration than SNAS... can authors provide more intuition on this part.

27 **R6:** The main reason is that only child architectures are sampled and optimized in DATA, which is faster than optimizing
28 the whole network in SNAS. As shown in Figure 2, DATA requires less training epochs than SNAS (100 vs 150).

29 Please kindly check and reassess our paper based on these requested experimental results. Thanks.

Model	Search	Child	Gap
DARTS	12.33	45.34	33.01
SNAS	11.46	9.33	2.13
DATA (M=7)	11.08	9.21	1.87

Model	Error (%)	Params (M)
DARTS(k=1)	3.00 ± 0.14	3.30
DARTS(k=2)	3.10 ± 0.12	4.00
DARTS(k=3)	2.95 ± 0.13	5.20
SNAS	2.85 ± 0.02	2.80
DATA(M=1)	2.95 ± 0.08	2.59
DATA(M=4)	2.70 ± 0.10	3.11
DATA(M=7)	2.59 ± 0.09	3.74



30 **To reviewer #2:** (Q: Why the operations used in each edge should be mutually exclusive)

31 **R:** As explained at line 125 in the paper, DATA is actually inspired by the fact that operations may NOT be mutually
32 exclusive but compatible, in contrast to SNAS and DARTS. In SNAS, any probability vector literally converges to an
33 one-hot vector, which means only one operation is selected. Though DARTS is capable of selecting the top- k strongest
34 operations, it always subjects to two limitations. First, the number of selected operations k is fixed on each edge, while
35 DATA can adaptively adjust the numbers on different edges. Second, DATA (gap=1.87) can seamlessly bridge the gap
36 between architectures in searching and validating, which obviously excels DARTS (gap=33.01), as shown in Table 1.

37 **To reviewer #2 and #3:** (Q: Comparison with a baseline that selects top- k operations during searching)

38 **R:** The results of selecting top- k operations using DARTS are reported in Table 2, which show that DARTS tends to
39 select superfluous while inefficient operations that may dramatically introduce more parameters. In contrast, our EGS
40 endows DATA with the capability of learning simple yet powerful operations, *e.g.*, the skip connection in Figure 1.

41 **To reviewer #3:** (Q: It would be interesting to know what other methods the authors tried... any convergence...)

42 **R:** Thanks for your interest. We have in the first place explored some naive methods of sampling multiple operations on
43 each edge, *i.e.*, selecting top- k operations. However, according to the issues as explained in our response to reviewer #2
44 and the results in Table 2, none of them works well in practice. For the convergence of DATA, we verify it empirically
45 in Figure 2 and 4 in this response. It can be seen that our method converges favorably well.

46 [C1] Bender, Gabriel, *et al.* "Understanding and simplifying one-shot architecture search." ICML. 2018.